

PCT

WORLD INTELLECTUAL PROPERTY ORGANIZATION  
International Bureau



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification <sup>7</sup> : <b>C12N</b>		<b>A2</b>	(11) International Publication Number: <b>WO 00/27994</b>
			(43) International Publication Date: 18 May 2000 (18.05.00)
(21) International Application Number: <b>PCT/US99/26923</b>			(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).
(22) International Filing Date: 12 November 1999 (12.11.99)			
(30) Priority Data: 60/108,279 12 November 1998 (12.11.98) US 60/128,606 8 April 1999 (08.04.99) US			
(71) Applicant: THE REGENTS OF THE UNIVERSITY OF CALIFORNIA [US/US]; 12th Floor, 1111 Franklin Street, Oakland, CA 94607 (US).			
(72) Inventors: STEPHENS, Richard MITCHELL, Wayne KALMAN, Sue DAVIS, Ronald			
(74) Agents: BASTIAN, Kevin, L. et al.; Townsend and Townsend and Crew LLP, 8th floor, Two Embarcadero Center, San Francisco, CA 94111-3834 (US).			Published <i>Without international search report and to be republished upon receipt of that report.</i>
(54) Title: CHLAMYDIA PNEUMONIAE GENOME SEQUENCE			
(57) Abstract  <p><i>C. pneumoniae</i> genome sequence and analysis of the encoded polypeptides and RNAs are provided. The <i>C. pneumoniae</i> gene nucleic acid compositions find use in identifying homologous or related proteins and the DNA sequences encoding such proteins; in producing compositions that modulate the expression or function of the protein; and in studying associated physiological pathways. In addition, modulation of the gene activity <i>in vivo</i> is used for prophylactic and therapeutic purposes, such as identification of cell type based on expression, and the like.</p>			

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

## CHLAMYDIA PNEUMONIAE GENOME SEQUENCE

## CROSS-REFERENCES TO RELATED APPLICATIONS

The present application is related to 60/128,606, filed April 8, 1999 and 60/108,279, filed November 12, 1998, which are incorporated herein by reference.

STATEMENT AS TO RIGHTS TO INVENTIONS MADE UNDER  
FEDERALLY SPONSORED RESEARCH AND DEVELOPMENT

## FIELD OF THE INVENTION

This invention relates to nucleic acids and polypeptides from *Chlamydia pneumoniae* and to their use in the diagnosis, prevention and treatment of diseases associated with *C. pneumoniae*.

## BACKGROUND OF THE INVENTION

*Chlamydiaceae* is a family of obligate intracellular parasite with a tropism for epithelial cells lining the mucus membranes. The bacteria have two morphologically distinct forms, "elementary body" and "reticulate body". The elementary body is the infectious form, and has a rigid cell wall, primarily of cross-linked outer membrane proteins. The reticulate body is the intracellular, metabolically active form. A unique developmental cycle between these two forms characterizes *Chlamydia* growth.

*C. pneumoniae* is a human respiratory pathogen that causes acute respiratory disease, and approximately 10% of community-acquired pneumonia. Antibody prevalence studies have shown that virtually everyone is infected with *C. pneumoniae* at some time, and that reinfection is common. In addition to respiratory disease, studies have shown an association of this organism with coronary artery disease. It has been demonstrated in atherosclerotic lesions of the aorta and coronary arteries by immunocytochemistry and by polymerase chain reaction (Kuo *et al.* (1993) *J Infect Dis* 167(4):841-849).

Recent reports have further demonstrated the presence of *C. pneumoniae* in the walls of abdominal aortic aneurysms (Juvonen *et al.* (1997) *J Vasc Surg* 25(3):499-505). Abdominal aortic aneurysms are frequently associated with atherosclerosis, and inflammation may be an important factor in aneurysmal dilatation.

5 *C. pneumoniae* may play a role in maintaining an inflammation and triggering the development of aortic aneurysms.

Muhlestein *et al.* (1996) JACC 27:1555-61, reported a differential incidence of *Chlamydia* species within the coronary artery wall of patients with atherosclerosis versus those with other forms of cardiovascular disease. The extremely high rate of possible infection in patients with symptomatic atherosclerotic disease compared to the very low rate in patients with normal coronary arteries or coronary artery disease from chronic transplant rejection provides evidence for a direct link between the atherosclerotic process and *Chlamydia* infection. Because a history of chlamydial infection is so prevalent in the population, the issue of causality remains. On a physiologic and pathologic level, abnormal interactions among endothelial cells, platelets, macrophages and lymphocytes may lead to a cascade of events resulting in acute endothelial damage, thrombosis and repair, chronically leading to the development of atheroma in blood vessels.

15 *C. pneumoniae* is related to other *Chlamydia* species, but the level of sequence similarity is relatively low. Very little is known about the biology of this organism, although it appears to be an important human pathogen. Allelic diversity and structural relationships between specific genes of Chlamydial species is described in Kaltenboeck *et al.* (1993) J Bacteriol 175(2):487-502; Gaydos *et al.* (1992) Infect Immun 60(12):5319-5323; Everett *et al.* (1997) Int J Syst Bacteriol 47(2):461-473; and Pudjiatmoko *et al.* (1997) Int J Syst Bacteriol 47(2):425-431.

A number of studies have been published describing methods for detection of *C. pneumoniae*, and for distinguishing between Chlamydial species. Such methods include PCR detection (Rasmussen *et al.* (1992) Mol Cell Probes 6(5):389-394; Holland *et al.* (1990) J Infect Dis 162(4):984-987); a simplified polymerase chain reaction-enzyme immunoassay (Wilson *et al.* (1996) J Appl Bacteriol 80(4):431-438); sequence determination and restriction endonuclease cleavage (Herrmann *et al.* (1996) J Clin Microbiol 34(8):1897-1902).

Antigenic and molecular analyses of different *C. pneumoniae* strains is described in Jantos *et al.* (1997) J Clin Microbiol 35(3):620-623. Some genes of *C. pneumoniae* have been isolated and sequenced. These include the Gro E operon (Kikuta *et al.* (1991) Infect Immun 59(12):4665-4669); the major outer membrane protein Perez *et*

5 *al.* (1991) Infect Immun 59(6):2195-2199; the DnaK protein homolog (Kornak *et al.*  
(1991) Infect Immun 59(2):721-725); as well as a number of ribosomal and other genes.

#### SUMMARY OF THE INVENTION

This invention provides the genomic sequence of *Chlamydia pneumoniae*.

15 The sequence information is useful for a variety of diagnostic and analytical methods.  
The genomic sequence may be embodied in a variety of media, including computer  
10 readable forms, or as a nucleic acid comprising a selected fragment of the sequence.  
Such fragments generally consist of an open reading frame, transcriptional or translational  
20 control elements, or fragments derived therefrom. Proteins encoded by the open reading  
frames are useful for diagnostic purposes, as well as for their enzymatic or structural  
activity.

#### DEFINITIONS

25 The term "amino acid" refers to naturally occurring and synthetic amino  
acids, as well as amino acid analogs and amino acid mimetics that function in a manner  
30 similar to the naturally occurring amino acids. Naturally occurring amino acids are those  
20 encoded by the genetic code, as well as those amino acids that are later modified, e.g.,  
hydroxyproline,  $\gamma$ -carboxyglutamate, and O-phosphoserine. Amino acid analogs refers to  
compounds that have the same basic chemical structure as a naturally occurring amino  
35 acid, i.e., an  $\alpha$  carbon that is bound to a hydrogen, a carboxyl group, an amino group, and  
an R group., e.g., homoserine, norleucine, methionine sulfoxide, methionine methyl  
25 sulfonium. Such analogs have modified R groups (e.g., norleucine) or modified peptide  
backbones, but retain the same basic chemical structure as a naturally occurring amino  
40 acid. Amino acid mimetics refers to chemical compounds that have a structure that is  
different from the general chemical structure of an amino acid, but that functions in a  
manner similar to a naturally occurring amino acid.

45 30 Amino acids may be referred to herein by either their commonly known  
three letter symbols or by the one-letter symbols recommended by the IUPAC-IUB  
Biochemical Nomenclature Commission. Nucleotides, likewise, may be referred to by  
50 their commonly accepted single-letter codes.

5 "Amplification" primers are oligonucleotides comprising either natural or analogue nucleotides that can serve as the basis for the amplification of a select nucleic acid sequence. They include, e.g., polymerase chain reaction primers and ligase chain reaction oligonucleotides.

10 5 "Antibody" refers to an immunoglobulin molecule able to bind to a specific epitope on an antigen. Antibodies can be a polyclonal mixture or monoclonal. Antibodies can be intact immunoglobulins derived from natural sources or from recombinant sources and can be immunoreactive portions of intact immunoglobulins. 15 Antibodies may exist in a variety of forms including, for example, Fv, Fab, and F(ab)<sub>2</sub>, as well as in single chains. Single-chain antibodies, in which genes for a heavy chain and a light chain are combined into a single coding sequence, may also be used.

20 An "antigen" is a molecule that is recognized and bound by an antibody, e.g., peptides, carbohydrates, organic molecules, or more complex molecules such as glycolipids and glycoproteins. The part of the antigen that is the target of antibody 25 binding is an antigenic determinant and a small functional group that corresponds to a single antigenic determinant is called a hapten.

"Biological sample" refers to any sample obtained from a living or dead organism. Examples of biological samples include biological fluids and tissue specimens. 30 Such biological samples can be prepared for analysis of the presence of *C. pneumoniae* nucleic acids, proteins, or antibodies specifically reactive with the proteins.

35 The term "*C. pneumoniae* gene" shall be intended to mean the open reading frame encoding specific *C. pneumoniae* polypeptides, as well as adjacent 5' and 3' non-coding nucleotide sequences involved in the regulation of expression, up to about 2 kb beyond the coding region, but possibly further in either direction. The gene may be 25 introduced into an appropriate vector for extrachromosomal maintenance or for integration into a host genome.

40 "Conservatively modified variants" applies to both amino acid and nucleic acid sequences. With respect to particular nucleic acid sequences, conservatively modified variants refers to those nucleic acids which encode identical or essentially 45 identical amino acid sequences, or where the nucleic acid does not encode an amino acid sequence, to essentially identical sequences. Specifically, degenerate codon substitutions may be achieved by generating sequences in which the third position of one or more 50 selected (or all) codons is substituted with mixed-base and/or deoxyinosine residues

(Batzner *et al.*, *Nucleic Acid Res.* 19:5081 (1991); Ohtsuka *et al.*, *J. Biol. Chem.* 260:2605-2608 (1985); Rossolini *et al.*, *Mol. Cell. Probes* 8:91-98 (1994)). Because of the degeneracy of the genetic code, a large number of functionally identical nucleic acids encode any given protein. For instance, the codons GCA, GCC, GCG and GCU all encode the amino acid alanine. Thus, at every position where an alanine is specified by a codon, the codon can be altered to any of the corresponding codons described without altering the encoded polypeptide. Such nucleic acid variations are "silent variations," which are one species of conservatively modified variations. Every nucleic acid sequence herein which encodes a polypeptide also describes every possible silent variation of the nucleic acid. One of skill will recognize that each codon in a nucleic acid (except AUG, which is ordinarily the only codon for methionine, and TGG, which is ordinarily the only codon for tryptophan) can be modified to yield a functionally identical molecule. Accordingly, each silent variation of a nucleic acid which encodes a polypeptide is implicit in each described sequence.

As to amino acid sequences, one of skill will recognize that individual substitutions, deletions or additions to a nucleic acid, peptide, polypeptide, or protein sequence which alters, adds or deletes a single amino acid or a small percentage of amino acids in the encoded sequence is a "conservatively modified variant" where the alteration results in the substitution of an amino acid with a chemically similar amino acid. Conservative substitution tables providing functionally similar amino acids are well known in the art. Such conservatively modified variants are in addition to and do not exclude polymorphic variants, interspecies homologs, and alleles of the invention.

The following groups each contain amino acids that are conservative substitutions for one another:

- 1) Alanine (A), Glycine (G);
- 2) Serine (S), Threonine (T);
- 3) Aspartic acid (D), Glutamic acid (E);
- 4) Asparagine (N), Glutamine (Q);
- 5) Cysteine (C), Methionine (M);
- 6) Arginine (R), Lysine (K), Histidine (H);
- 7) Isoleucine (I), Leucine (L), Valine (V); and
- 8) Phenylalanine (F), Tyrosine (Y), Tryptophan (W).

see, e.g., Creighton, *Proteins* (1984)).

5                   The terms "identical" or percent "identity," in the context of two or more  
nucleic acids or polypeptide sequences, refer to two or more sequences or subsequences  
that are the same or have a specified percentage of amino acid residues or nucleotides that  
10                   are the same, when compared and aligned for maximum correspondence over a  
5                   comparison window, as measured using one of the following sequence comparison  
algorithms or by manual alignment and visual inspection. This definition also refers to  
the complement of a test sequence, which has a designated percent sequence or  
15                   subsequence complementarity when the test sequence has a designated or substantial  
identity to a reference sequence. For example, a designated amino acid percent identity  
10                   of 95% refers to sequences or subsequences that have at least about 95% amino acid  
identity when aligned for maximum correspondence over a comparison window as  
20                   measured using one of the following sequence comparison algorithms or by manual  
alignment and visual inspection. Such sequences would then be said to have substantial  
identity, or to be substantially identical to each other. Preferably, sequences have at least  
25                   15                   about 70% identity, more preferably 80% identity, more preferably 90-95% identity and  
above. Preferably, the percent identity exists over a region of the sequence that is at least  
about 25 amino acids in length, more preferably over a region that is 50-100 amino acids  
in length.

30                   When percentage of sequence identity is used in reference to proteins or  
20                   peptides, it is recognized that residue positions that are not identical often differ by  
conservative amino acid substitutions, where amino acids residues are substituted for  
35                   other amino acid residues with similar chemical properties (e.g., charge or  
hydrophobicity) and therefore do not change the functional properties of the molecule.  
Where sequences differ in conservative substitutions, the percent sequence identity may  
25                   be adjusted upwards to correct for the conservative nature of the substitution. Means for  
40                   making this adjustment are well known to those of skill in the art. Typically this involves  
scoring a conservative substitution as a partial rather than a full mismatch, thereby  
increasing the percentage sequence identity. Thus, for example, where an identical amino  
45                   acid is given a score of 1 and a non-conservative substitution is given a score of zero, a  
30                   conservative substitution is given a score between zero and 1. The scoring of  
conservative substitutions is calculated according to, e.g., the algorithm of Meyers &  
Miller, *Computer Applic. Biol. Sci.* 4:11-17 (1988) e.g., as implemented in the program  
50                   PC/GENE (Intelligenetics, Mountain View, California, USA)..



5 For sequence comparison, typically one sequence acts as a reference  
sequence, to which test sequences are compared. When using a sequence comparison  
algorithm, test and reference sequences are entered into a computer, subsequence  
10 coordinates are designated, if necessary, and sequence algorithm program parameters are  
5 designated. Default program parameters can be used, or alternative parameters can be  
designated. The sequence comparison algorithm then calculates the percent sequence  
identity for the test sequence(s) relative to the reference sequence, based on the  
15 designated or default program parameters.

A comparison window includes reference to a segment of any one of the  
10 number of contiguous positions selected from the group consisting of from 25 to 600,  
usually about 50 to about 200, more usually about 100 to about 150 in which a sequence  
20 may be compared to a reference sequence of the same number of contiguous positions  
after the two sequences are optimally aligned. Methods of alignment of sequences for  
comparison are well-known in the art. Optimal alignment of sequences for comparison  
25 can be conducted, e.g., by the local homology algorithm of Smith & Waterman, *Adv.  
Appl. Math.* 2:482 (1981), by the homology alignment algorithm of Needleman &  
Wunsch, *J. Mol. Biol.* 48:443 (1970), by the search for similarity method of Pearson &  
Lipman, *Proc. Nat'l. Acad. Sci. USA* 85:2444 (1988), by computerized implementations  
30 of these algorithms (GAP, BESTFIT, FASTA, and TFASTA in the Wisconsin Genetics  
20 Software Package, Genetics Computer Group, 575 Science Dr., Madison, WI), or by  
manual alignment and visual inspection (*see, e.g., Ausubel et al., supra*).

One example of a useful algorithm is PILEUP. PILEUP creates a multiple  
35 sequence alignment from a group of related sequences using progressive, pairwise  
alignments to show relationship and percent sequence identity. It also plots a tree or  
25 dendrogram showing the clustering relationships used to create the alignment. PILEUP  
40 uses a simplification of the progressive alignment method of Feng & Doolittle, *J. Mol.  
Evol.* 35:351-360 (1987). The method used is similar to the method described by Higgins  
& Sharp, *CABIOS* 5:151-153 (1989). The program can align up to 300 sequences, each  
45 of a maximum length of 5,000 nucleotides or amino acids. The multiple alignment  
30 procedure begins with the pairwise alignment of the two most similar sequences,  
producing a cluster of two aligned sequences. This cluster is then aligned to the next  
most related sequence or cluster of aligned sequences. Two clusters of sequences are  
50 aligned by a simple extension of the pairwise alignment of two individual sequences. The

5 final alignment is achieved by a series of progressive, pairwise alignments. The program  
is run by designating specific sequences and their amino acid or nucleotide coordinates  
for regions of sequence comparison and by designating the program parameters. Using  
10 PILEUP, a reference sequence is compared to other test sequences to determine the  
5 percent sequence identity relationship using the following parameters: default gap weight  
(3.00), default gap length weight (0.10), and weighted end gaps. PILEUP can be obtained  
from the GCG sequence analysis software package, e.g., version 7.0 (Devereaux *et al.*,  
15 *Nuc. Acids Res.* 12:387-395 (1984).

Another example of algorithm that is suitable for determining percent  
10 sequence identity (i.e., substantial similarity or identity) is the BLAST algorithm, which  
is described in Altschul *et al.*, *J. Mol. Biol.* 215:403-410 (1990). Software for performing  
20 BLAST analyses is publicly available through the National Center for Biotechnology  
Information (<http://www.ncbi.nlm.nih.gov/>). This algorithm involves first identifying  
high scoring sequence pairs (HSPs) by identifying short words of length W in the query  
25 sequence, which either match or satisfy some positive-valued threshold score T when  
aligned with a word of the same length in a database sequence. T is referred to as the  
neighborhood word score threshold (Altschul *et al, supra*). These initial neighborhood  
word hits act as seeds for initiating searches to find longer HSPs containing them. The  
30 word hits are then extended in both directions along each sequence for as far as the  
20 cumulative alignment score can be increased. Cumulative scores are calculated using, for  
nucleotide sequences, the parameters M (reward score for a pair of matching residues;  
always > 0) and N (penalty score for mismatching residues, always < 0). For amino acid  
35 sequences, a scoring matrix is used to calculate the cumulative score. Extension of the  
word hits in each direction are halted when: the cumulative alignment score falls off by  
25 the quantity X from its maximum achieved value; the cumulative score goes to zero or  
below, due to the accumulation of one or more negative-scoring residue alignments; or  
40 the end of either sequence is reached. The BLAST algorithm parameters W, T, and X  
determine the sensitivity and speed of the alignment. The BLASTN program (for  
nucleotide sequences) uses as defaults a wordlength (W) of 11, an expectation (E) of 10,  
45 M=5, N=4, and a comparison of both strands. For amino acid sequences, the BLASTP  
program uses as default parameters a wordlength (W) of 3, an expectation (E) of 10, and  
the BLOSUM62 scoring matrix (*see* Henikoff & Henikoff, *Proc. Natl. Acad. Sci. USA*  
50 89:10915 (1989)).

5 The BLAST algorithm also performs a statistical analysis of the similarity  
between two sequences (see, e.g., Karlin & Altschul, *Proc. Nat'l. Acad. Sci. USA*  
90:5873-5787 (1993)). One measure of similarity provided by the BLAST algorithm is  
10 the smallest sum probability (P(N)), which provides an indication of the probability by  
5 which a match between two nucleotide or amino acid sequences would occur by chance.  
For example, a nucleic acid is considered similar to a reference sequence if the smallest  
sum probability in a comparison of the test nucleic acid to the reference nucleic acid is  
15 less than about 0.1, more preferably less than about 0.01, and most preferably less than  
about 0.001.

10 An indication that two nucleic acid sequences or polypeptides are  
substantially identical is that the polypeptide encoded by the first nucleic acid is  
20 immunologically cross reactive with the antibodies raised against the polypeptide  
encoded by the second nucleic acid, as described below. Thus, a polypeptide is typically  
substantially identical to a second polypeptide, for example, where the two peptides differ  
25 only by conservative substitutions. Another indication that two nucleic acid sequences  
are substantially identical is that the two molecules or their complements hybridize to  
each other under stringent conditions, as described below.

30 Another indication that polynucleotide sequences are substantially  
identical is if two molecules hybridize to each other under stringent conditions. Stringent  
20 conditions are sequence dependent and will be different in different circumstances.  
Generally, stringent conditions are selected to be about 5°C lower than the thermal  
melting point (T<sub>m</sub>) for the specific sequence at a defined ionic strength and pH. The T<sub>m</sub>  
35 is the temperature (under defined ionic strength and pH) at which 50% of the target  
sequence hybridizes to a perfectly matched probe. Typically stringent conditions for a  
25 Southern blot protocol involve hybridizing in a buffer comprising 5x SSC, 1% SDS at  
40 65°C or hybridizing in a buffer containing 5x SSC and 1% SDS at 42°C and washing at  
65°C with a 0.2x SSC, 0.1% SDS wash.

45 A "label" is a composition detectable by spectroscopic, photochemical,  
biochemical, immunochemical, or chemical means. For example, useful labels include  
30 <sup>32</sup>P, fluorescent dyes, electron-dense reagents, enzymes (e.g., as commonly used in an  
ELISA), biotin, dioxigenin, or haptens and proteins for which antisera or monoclonal  
antibodies are available.

5                   The term "nucleic acid" refers to deoxyribonucleotides or ribonucleotides  
and polymers thereof in either single- or double-stranded form. The term encompasses  
nucleic acids containing known nucleotide analogs or modified backbone residues or  
10 linkages, which are synthetic, naturally occurring, and non-naturally occurring, which  
5 have similar binding properties as the reference nucleic acid, and which are metabolized  
in a manner similar to the reference nucleotides. Examples of such analogs include,  
without limitation, phosphorothioates, phosphoramidates, methyl phosphonates,  
15 chiral-methyl phosphonates, 2-O-methyl ribonucleotides, peptide-nucleic acids (PNAs).

                  Unless otherwise indicated, a particular nucleic acid sequence also  
10 implicitly encompasses conservatively modified variants thereof (e.g., degenerate codon  
substitutions) and complementary sequences, as well as the sequence explicitly indicated.  
20 The term nucleic acid is used interchangeably with gene, cDNA, mRNA, oligonucleotide,  
and polynucleotide.

                  As used herein a "nucleic acid probe or oligonucleotide" is defined as a  
25 nucleic acid capable of binding to a target nucleic acid of complementary sequence  
through one or more types of chemical bonds, usually through complementary base  
pairing, usually through hydrogen bond formation. As used herein, a probe may include  
30 natural (i.e., A, G, C, or T) or modified bases (7-deazaguanosine, inosine, etc.). In  
addition, the bases in a probe may be joined by a linkage other than a phosphodiester  
20 bond, so long as it does not interfere with hybridization. Thus, for example, probes may  
be peptide nucleic acids in which the constituent bases are joined by peptide bonds rather  
than phosphodiester linkages. It will be understood by one of skill in the art that probes  
35 may bind target sequences lacking complete complementarity with the probe sequence  
depending upon the stringency of the hybridization conditions. The probes are preferably  
25 directly labeled as with isotopes, chromophores, lumiphores, chromogens, or indirectly  
40 labeled such as with biotin to which a streptavidin complex may later bind. By assaying  
for the presence or absence of the probe, one can detect the presence or absence of the  
select sequence or subsequence.

                  A labeled nucleic acid probe or oligonucleotide is one that is bound, either  
45 covalently, through a linker, or through ionic, van der Waals or hydrogen bonds to a label  
30 such that the presence of the probe may be detected by detecting the presence of the label  
bound to the probe.

5 "Pharmaceutically acceptable" means a material that is not biologically or otherwise undesirable, i.e., the material can be administered to an individual along with a *Chlamydia* antigen without causing any undesirable biological effects or interacting in a deleterious manner with any of the other components of the pharmaceutical composition.

10 5 The terms "polypeptide," "peptide" and "protein" are used interchangeably herein to refer to a polymer of amino acid residues. The terms apply to amino acid polymers in which one or more amino acid residue is an analog or mimetic of a corresponding naturally occurring amino acid, as well as to naturally occurring amino acid polymers.

15 10 The phrase "specifically or selectively hybridizing to," refers to hybridization between a probe and a target sequence in which the probe binds substantially only to the target sequence, forming a hybridization complex, when the target is in a heterogeneous mixture of polynucleotides and other compounds. Such hybridization is determinative of the presence of the target sequence. Although the probe  
20 15 may bind other unrelated sequences, at least 90%, preferably 95% or more of the hybridization complexes formed are with the target sequence.

25 30 The term "recombinant" when used with reference to a cell, or nucleic acid, or vector, indicates that the cell, or nucleic acid, or vector, has been modified by the introduction of a heterologous nucleic acid or the alteration of a native nucleic acid, or  
35 20 that the cell is derived from a cell so modified. Thus, for example, recombinant cells express genes that are not found within the native (non-recombinant) form of the cell or express native genes that are otherwise abnormally expressed, under expressed or not expressed at all.

40 25 The phrase "specifically immunoreactive with", when referring to a protein or peptide, refers to a binding reaction between the protein and an antibody which is determinative of the presence of the protein in the presence of a heterogeneous population of proteins and other compounds. Thus, under designated immunoassay conditions, the specified antibodies bind to a particular protein and do not bind in a significant amount to other proteins present in the sample. Specific binding to an antibody under such  
45 30 conditions may require an antibody that is selected for its specificity for a particular protein. A variety of immunoassay formats may be used to select antibodies specifically immunoreactive with a particular protein and are described in detail below.

5 The phrase "substantially pure" or "isolated" when referring to a  
subcellular components of the *Chlamydia* organism. Typically, a monomeric protein is  
substantially pure when at least about 85% or more of a sample exhibits a single  
10 polypeptide backbone. Minor variants or chemical modifications may typically share the  
5 same polypeptide sequence. Depending on the purification procedure, purities of 85%,  
and preferably over 95% pure are possible. Protein purity or homogeneity may be  
15 indicated by a number of means well known in the art, such as polyacrylamide gel  
electrophoresis of a protein sample, followed by visualizing a single polypeptide band on  
10 a polyacrylamide gel upon silver staining. For certain purposes high resolution will be  
needed and HPLC or a similar means for purification utilized.

#### DETAILED DESCRIPTION

The present invention provides the nucleotide sequence of the *C.*  
25 *pneumoniae* genome SEQ ID NO: 1 or a representative fragment thereof, in a form which  
can be readily used, analyzed, and interpreted by a skilled artisan. As used herein, a  
"representative fragment" of the nucleotide sequence depicted in SEQ ID NO: 1 refers to  
any portion which is not presently represented within a publicly available database.  
30 Preferred representative fragments of the present invention are open reading frames,  
expression modulating fragments, uptake modulating fragments, and fragments which can  
20 be used to diagnose the presence of *C. pneumoniae* in sample. Using the information  
provided in the present application, together with routine cloning and sequencing  
35 methods, one of ordinary skill in the art will be able to clone and sequence all  
"representative fragments" of interest including open reading frames (ORFs) encoding a  
25 large variety of *C. pneumoniae* proteins. A non-limiting identification of such preferred  
representative fragments is provided in Tables 2 and 3.

#### Diagnostic use of *C. pneumoniae* nucleic acids

##### Hybridization-based assays

45 Using the nucleic acids disclosed here, one of skill can design nucleic acid  
30 hybridization-based assays for the detection of *C. pneumoniae*. Any of a number of well  
known techniques for the specific detection of target nucleic acids can be used.  
50 Exemplary hybridization-based assays include, but are not limited to, traditional "direct

probe" methods such as Southern Blots, dot blots, *in situ* hybridization (e.g., FISH), PCR, and the like. The methods can be used in a wide variety of formats including, but not limited to substrate- (e.g. membrane or glass) bound methods or array-based approaches as described below. As noted above, this invention also embraces methods for detecting the presence of *Chlamydia* DNA or RNA in biological samples. These sequences can be used to detect *Chlamydia* in biological samples from patients suspected of being infected. A variety of methods of specific DNA and RNA measurement using nucleic acid hybridization techniques are known to those of skill in the art (see Sambrook *et al.*, *supra*).

*In situ* hybridization assays are well known (e.g., Angerer (1987) *Meth. Enzymol* 152: 649). Generally, *in situ* hybridization comprises the following major steps: (1) fixation of tissue or biological structure to be analyzed; (2) prehybridization treatment of the biological structure to increase accessibility of target DNA, and to reduce nonspecific binding; (3) hybridization of the mixture of nucleic acids to the nucleic acid in the biological structure or tissue; (4) post-hybridization washes to remove nucleic acid fragments not bound in the hybridization and (5) detection of the hybridized nucleic acid fragments. The reagent used in each of these steps and the conditions for use vary depending on the particular application.

In a typical *in situ* hybridization assay, cells are fixed to a solid support, typically a glass slide. If a nucleic acid is to be probed, the cells are typically denatured with heat or alkali. The cells are then contacted with a hybridization solution at a moderate temperature to permit annealing of labeled probes specific to the nucleic acid sequence encoding the protein. The targets (e.g., cells) are then typically washed at a predetermined stringency or at an increasing stringency until an appropriate signal to noise ratio is obtained.

The nucleic acids of this invention are particularly well suited to array-based hybridization formats. Arrays are a multiplicity of different "probe" or "target" nucleic acids (or other compounds) attached to one or more surfaces (e.g., solid, membrane, or gel). In a preferred embodiment, the multiplicity of nucleic acids (or other moieties) is attached to a single contiguous surface or to a multiplicity of surfaces juxtaposed to each other.

In an array format a large number of different hybridization reactions can be run essentially "in parallel." This provides rapid, essentially simultaneous, evaluation

5 of a number of hybridizations in a single "experiment". Methods of performing hybridization reactions in array based formats are well known to those of skill in the art (see, e.g., Pastinen (1997) *Genome Res.* 7: 606-614; Jackson (1996) *Nature Biotechnology* 14:1685; Chee (1995) *Science* 274: 610; WO 96/17958.

10 5 Arrays, particularly nucleic acid arrays can be produced according to a wide variety of methods well known to those of skill in the art. For example, in a simple embodiment, "low density" arrays can simply be produced by spotting (e.g. by hand using a pipette) different nucleic acids at different locations on a solid support (e.g. a glass surface, a membrane, etc.).

15 10 This simple spotting, approach has been automated to produce high density spotted arrays (see, e.g., U.S. Patent No: 5,807,522). This patent describes the use of an automated systems that taps a microcapillary against a surface to deposit a small volume of a biological sample. The process is repeated to generate high density arrays. Arrays can also be produced using oligonucleotide synthesis technology. Thus, for  
20 15 example, U.S. Patent No. 5,143,854 and PCT patent publication Nos. WO 90/15070 and 92/10092 teach the use of light-directed combinatorial synthesis of high density oligonucleotide arrays.

30 20 Many methods for immobilizing nucleic acids on a variety of solid surfaces are known in the art. A wide variety of organic and inorganic polymers, as well as other materials, both natural and synthetic, can be employed as the material for the solid surface. Illustrative solid surfaces include, e.g., nitrocellulose, nylon, glass, quartz, diazotized membranes (paper or nylon), silicones, polyformaldehyde, cellulose, and  
35 25 cellulose acetate. In addition, plastics such as polyethylene, polypropylene, polystyrene, and the like can be used. Other materials which may be employed include paper, ceramics, metals, metalloids, semiconductive materials, cermets or the like. In addition, substances that form gels can be used. Such materials include, e.g., proteins (e.g.,  
40 30 gelatins), lipopolysaccharides, silicates, agarose and polyacrylamides. Where the solid surface is porous, various pore sizes may be employed depending upon the nature of the system.

45 30 In preparing the surface, a plurality of different materials may be employed, particularly as laminates, to obtain various properties. For example, proteins (e.g., bovine serum albumin) or mixtures of macromolecules (e.g., Denhardt's solution)  
50 35 can be employed to avoid non-specific binding, simplify covalent conjugation, enhance



5 signal detection or the like. If covalent bonding between a compound and the surface is  
desired, the surface will usually be polyfunctional or be capable of being  
polyfunctionalized. Functional groups which may be present on the surface and used for  
10 linking can include carboxylic acids, aldehydes, amino groups, cyano groups, ethylenic  
5 groups, hydroxyl groups, mercapto groups and the like. The manner of linking a wide  
variety of compounds to various surfaces is well known and is amply illustrated in the  
literature.

15 For example, methods for immobilizing nucleic acids by introduction of  
various functional groups to the molecules is known (*see, e.g.*, Bischoff (1987) *Anal.*  
10 *Biochem.*, 164: 336-344; Kremsky (1987) *Nucl. Acids Res.* 15: 2891-2910). Modified  
nucleotides can be placed on the target using PCR primers containing the modified  
20 nucleotide, or by enzymatic end labeling with modified nucleotides. Use of glass or  
membrane supports (*e.g.*, nitrocellulose, nylon, polypropylene) for the nucleic acid arrays  
of the invention is advantageous because of well developed technology employing  
25 manual and robotic methods of arraying targets at relatively high element densities. Such  
membranes are generally available and protocols and equipment for hybridization to  
membranes is well known.

30 Target elements of various sizes, ranging from 1 mm diameter down to 1  
 $\mu\text{m}$  can be used. Smaller target elements containing low amounts of concentrated, fixed  
20 probe DNA are used for high complexity comparative hybridizations since the total  
amount of sample available for binding to each target element will be limited. Thus it is  
advantageous to have small array target elements that contain a small amount of  
35 concentrated probe DNA so that the signal that is obtained is highly localized and bright.  
Such small array target elements are typically used in arrays with densities greater than  
25  $10^4/\text{cm}^2$ . Relatively simple approaches capable of quantitative fluorescent imaging of 1  
40  $\text{cm}^2$  areas have been described that permit acquisition of data from a large number of  
target elements in a single image (*see, e.g.*, Wittrup (1994) *Cytometry* 16:206-213).

45 If fluorescently labeled nucleic acid samples are used, arrays on solid  
surface substrates with much lower fluorescence than membranes, such as glass, quartz,  
30 or small beads, can achieve much better sensitivity. Substrates such as glass or fused  
silica are advantageous in that they provide a very low fluorescence substrate, and a  
highly efficient hybridization environment. Covalent attachment of the target nucleic  
50 acids to glass or synthetic fused silica can be accomplished according to a number of

known techniques (described above). Nucleic acids can be conveniently coupled to glass using commercially available reagents. For instance, materials for preparation of silanized glass with a number of functional groups are commercially available or can be prepared using standard techniques (see, e.g., Gait (1984) *Oligonucleotide Synthesis: A Practical Approach*, IRL Press, Wash., D.C.). Quartz cover slips, which have at least 10-fold lower autofluorescence than glass, can also be silanized.

Alternatively, probes can also be immobilized on commercially available coated beads or other surfaces. For instance, biotin end-labeled nucleic acids can be bound to commercially available avidin-coated beads. Streptavidin or anti-digoxigenin antibody can also be attached to silanized glass slides by protein-mediated coupling using e.g., protein A following standard protocols (see, e.g., Smith (1992) *Science* 258: 1122-1126). Biotin or digoxigenin end-labeled nucleic acids can be prepared according to standard techniques. Hybridization to nucleic acids attached to beads is accomplished by suspending them in the hybridization mix, and then depositing them on the glass substrate for analysis after washing. Alternatively, paramagnetic particles, such as ferric oxide particles, with or without avidin coating, can be used.

A variety of other nucleic acid hybridization formats are known to those skilled in the art. For example, common formats include sandwich assays and competition or displacement assays. Hybridization techniques are generally described in Hames and Higgins (1985) *Nucleic Acid Hybridization, A Practical Approach*, IRL Press; Gall and Pardue (1969) *Proc. Natl. Acad. Sci. USA* 63: 378-383; and John *et al.* (1969) *Nature* 223: 582-587.

Sandwich assays are commercially useful hybridization assays for detecting or isolating nucleic acid sequences. Such assays utilize a "capture" nucleic acid covalently immobilized to a solid support and a labeled "signal" nucleic acid in solution. The sample will provide the target nucleic acid. The "capture" nucleic acid and "signal" nucleic acid probe hybridize with the target nucleic acid to form a "sandwich" hybridization complex. To be most effective, the signal nucleic acid should not hybridize with the capture nucleic acid.

Detection of a hybridization complex may require the binding of a signal generating complex to a duplex of target and probe polynucleotides or nucleic acids. Typically, such binding occurs through ligand and anti-ligand interactions as between a ligand-conjugated probe and an anti-ligand conjugated with a signal.

5 The sensitivity of the hybridization assays may be enhanced through use of  
a nucleic acid amplification system that multiplies the target nucleic acid being detected.  
Examples of such systems include the polymerase chain reaction (PCR) system and the  
ligase chain reaction (LCR) system. Other methods recently described in the art are the  
10 nucleic acid sequence based amplification (NASBAO, Cangene, Mississauga, Ontario)  
and Q Beta Replicase systems.

Nucleic acid hybridization simply involves providing a denatured probe  
and target nucleic acid under conditions where the probe and its complementary target  
15 can form stable hybrid duplexes through complementary base pairing. The nucleic acids  
that do not form hybrid duplexes are then washed away leaving the hybridized nucleic  
acids to be detected, typically through detection of an attached detectable label. It is  
generally recognized that nucleic acids are denatured by increasing the temperature or  
decreasing the salt concentration of the buffer containing the nucleic acids, or in the  
20 addition of chemical agents, or the raising of the pH. Under low stringency conditions  
(e.g., low temperature and/or high salt and/or high target concentration) hybrid duplexes  
(e.g., DNA:DNA, RNA:RNA, or RNA:DNA) will form even where the annealed  
sequences are not perfectly complementary. Thus specificity of hybridization is reduced  
at lower stringency. Conversely, at higher stringency (e.g., higher temperature or lower  
25 salt) successful hybridization requires fewer mismatches.

One of skill in the art will appreciate that hybridization conditions may be  
selected to provide any degree of stringency. In a preferred embodiment, hybridization is  
performed at low stringency to ensure hybridization and then subsequent washes are  
performed at higher stringency to eliminate mismatched hybrid duplexes. Successive  
washes may be performed at increasingly higher stringency (e.g., down to as low as 0.25  
25 X SSPE-T at 37°C to 70°C) until a desired level of hybridization specificity is obtained.  
Stringency can also be increased by addition of agents such as formamide. Hybridization  
specificity may be evaluated by comparison of hybridization to the test probes with  
hybridization to the various controls that can be present.

In general, there is a tradeoff between hybridization specificity  
45 (stringency) and signal intensity. Thus, in a preferred embodiment, the wash is performed  
at the highest stringency that produces consistent results and that provides a signal  
intensity greater than approximately 10% of the background intensity. Thus, in a  
preferred embodiment, the hybridized array may be washed at successively higher  
50

5 stringency solutions and read between each wash. Analysis of the data sets thus produced will reveal a wash stringency above which the hybridization pattern is not appreciably altered and which provides adequate signal for the particular probes of interest.

10 Methods of optimizing hybridization conditions are well known to those of skill in the art (see, e.g., Tijssen (1993) *Laboratory Techniques in Biochemistry and Molecular Biology*, Vol. 24: *Hybridization With Nucleic Acid Probes*, Elsevier, N.Y.).

Labeling and detection of nucleic acids.

15 In a preferred embodiment, the hybridized nucleic acids are detected by detecting one or more labels attached to the sample or probe nucleic acids. The labels may be incorporated by any of a number of means well known to those of skill in the art. Means of attaching labels to nucleic acids include, for example nick translation or end-labeling (e.g. with a labeled RNA) by kinasing of the nucleic acid and subsequent attachment (ligation) of a nucleic acid linker joining the sample nucleic acid to a label (e.g., a fluorophore). A wide variety of linkers for the attachment of labels to nucleic acids are also known. In addition, intercalating dyes and fluorescent nucleotides can also be used.

20 Detectable labels suitable for use in the present invention include any composition detectable by spectroscopic, photochemical, biochemical, immunochemical, electrical, optical or chemical means. Useful labels in the present invention include biotin for staining with labeled streptavidin conjugate, magnetic beads (e.g., Dynabeads™), fluorescent dyes (e.g., fluorescein, texas red, rhodamine, green fluorescent protein, and the like, see, e.g., Molecular Probes, Eugene, Oregon, USA), radiolabels (e.g.,  $^3\text{H}$ ,  $^{125}\text{I}$ ,  $^{35}\text{S}$ ,  $^{14}\text{C}$ , or  $^{32}\text{P}$ ), enzymes (e.g., horse radish peroxidase, alkaline phosphatase and others commonly used in an ELISA), and colorimetric labels such as colloidal gold (e.g., gold particles in the 40 -80 nm diameter size range scatter green light with high efficiency) or colored glass or plastic (e.g., polystyrene, polypropylene, latex, etc.) beads. Patents teaching the use of such labels include U.S. Patent Nos. 3,817,837; 3,850,752; 3,939,350; 3,996,345; 4,277,437; 4,275,149; and 4,366,241.

25 A fluorescent label is preferred because it provides a very strong signal with low background. It is also optically detectable at high resolution and sensitivity through a quick scanning procedure. The nucleic acid samples can all be labeled with a single label, e.g., a single fluorescent label. Alternatively, in another embodiment, different nucleic acid samples can be simultaneously hybridized where each nucleic acid

5 sample has a different label. For instance, one target could have a green fluorescent label and a second target could have a red fluorescent label. The scanning step will distinguish  
10 cites of binding of the red label from those binding the green fluorescent label. Each nucleic acid sample (target nucleic acid) can be analyzed independently from one another.

15 — Suitable chromogens which can be employed include those molecules and compounds which absorb light in a distinctive range of wavelengths so that a color can be observed or, alternatively, which emit light when irradiated with radiation of a particular  
20 wave length or wave length range, e.g., fluorescers.

Desirably, fluorescers should absorb light above about 300 nm, preferably  
25 about 350 nm, and more preferably above about 400 nm, usually emitting at wavelengths greater than about 10 nm higher than the wavelength of the light absorbed. It should be noted that the absorption and emission characteristics of the bound dye can differ from  
30 the unbound dye. Therefore, when referring to the various wavelength ranges and characteristics of the dyes, it is intended to indicate the dyes as employed and not the dye  
35 which is unconjugated and characterized in an arbitrary solvent.

Fluorescers are generally preferred because by irradiating a fluorescer with  
40 light, one can obtain a plurality of emissions. Thus, a single label can provide for a plurality of measurable events.

Detectable signal can also be provided by chemiluminescent and  
45 bioluminescent sources. Chemiluminescent sources include a compound which becomes electronically excited by a chemical reaction and can then emit light which serves as the detectable signal or donates energy to a fluorescent acceptor. Alternatively, luciferins can  
50 be used in conjunction with luciferase or lucigenins to provide bioluminescence. Spin labels are provided by reporter molecules with an unpaired electron spin which can  
55 be detected by electron spin resonance (ESR) spectroscopy. Exemplary spin labels include organic free radicals, transitional metal complexes, particularly vanadium, copper, iron, and manganese, and the like. Exemplary spin labels include nitroxide free radicals.

The label may be added to the target (sample) nucleic acid(s) prior to, or  
60 after the hybridization. So called "direct labels" are detectable labels that are directly attached to or incorporated into the target (sample) nucleic acid prior to hybridization. In contrast, so called "indirect labels" are joined to the hybrid duplex after hybridization.  
65 Often, the indirect label is attached to a binding moiety that has been attached to the

target nucleic acid prior to the hybridization. Thus, for example, the target nucleic acid may be biotinylated before the hybridization. After hybridization, an avidin-conjugated fluorophore will bind the biotin bearing hybrid duplexes providing a label that is easily detected. For a detailed review of methods of labeling nucleic acids and detecting labeled hybridized nucleic acids see *Laboratory Techniques in Biochemistry and Molecular Biology, Vol. 24: Hybridization With Nucleic Acid Probes*, P. Tijssen, ed. Elsevier, N.Y., (1993)).

Fluorescent labels are easily added during an *in vitro* transcription reaction. Thus, for example, fluorescein labeled UTP and CTP can be incorporated into the RNA produced in an *in vitro* transcription.

The labels can be attached directly or through a linker moiety. In general, the site of label or linker-label attachment is not limited to any specific position. For example, a label may be attached to a nucleoside, nucleotide, or analogue thereof at any position that does not interfere with detection or hybridization as desired. For example, certain Label-ON Reagents from Clontech (Palo Alto, CA) provide for labeling interspersed throughout the phosphate backbone of an oligonucleotide and for terminal labeling at the 3' and 5' ends. As shown for example herein, labels can be attached at positions on the ribose ring or the ribose can be modified and even eliminated as desired. The base moieties of useful labeling reagents can include those that are naturally occurring or modified in a manner that does not interfere with the purpose to which they are put. Modified bases include but are not limited to 7-deaza A and G, 7-deaza-8-aza A and G, and other heterocyclic moieties.

It will be recognized that fluorescent labels are not to be limited to single species organic molecules, but include inorganic molecules, multi-molecular mixtures of organic and/or inorganic molecules, crystals, heteropolymers, and the like. Thus, for example, CdSe-CdS core-shell nanocrystals enclosed in a silica shell can be easily derivatized for coupling to a biological molecule (Bruchez *et al.* (1998) *Science*, 281: 2013-2016). Similarly, highly fluorescent quantum dots (zinc sulfide-capped cadmium selenide) have been covalently coupled to biomolecules for use in ultrasensitive biological detection (Warren and Nie (1998) *Science*, 281: 2016-2018).

#### Amplification-based assays.

In another embodiment, amplification-based assays can be used to detect nucleic acids. In such amplification-based assays, the nucleic acid sequences act as a

5 template in an amplification reaction (e.g. Polymerase Chain Reaction (PCR). Detailed protocols for quantitative PCR are provided in Innis *et al.* (1990) *PCR Protocols, A Guide to Methods and Applications*, Academic Press, Inc. N.Y.).

10 Other suitable amplification methods include, but are not limited to ligase chain reaction (LCR) (see Wu and Wallace (1989) *Genomics* 4: 560, Landegren *et al.* (1988) *Science* 241: 1077, and Barringer *et al.* (1990) *Gene* 89: 117, transcription amplification (Kwoh *et al.* (1989) *Proc. Natl. Acad. Sci. USA* 86: 1173), and self-sustained sequence replication (Guatelli *et al.* (1990) *Proc. Nat. Acad. Sci. USA* 87: 1874).

#### 10 Detection of *C. pneumoniae* gene expression

20 The nucleic acids of the invention can also be used to *C. pneumoniae* detect gene transcripts. Methods of detecting and/or quantifying gene transcripts using nucleic acid hybridization techniques are known to those of skill in the art (see Sambrook *et al. supra*). For example, a Northern transfer may be used for the detection of the  
25 15 desired mRNA directly. In brief, the mRNA is isolated from a given cell sample using, for example, an acid guanidinium-phenol-chloroform extraction method. The mRNA is then electrophoresed to separate the mRNA species and the mRNA is transferred from the gel to a nitrocellulose membrane. As with the Southern blots, labeled probes are used to  
30 identify and/or quantify the target mRNA.

20 In another preferred embodiment, the gene transcript can be measured using amplification (e.g. PCR) based methods as described above for directly assessing  
35 copy number of the target sequences.

#### Expression of *C. pneumoniae* proteins

40 25 The nucleic acids disclosed here can be used for recombinant expression of the proteins. In these methods, the nucleic acids encoding the proteins of interest are introduced into suitable host cells, followed by induction of the cells to produce large amounts of the protein. The invention relies on routine techniques in the field of recombinant genetics, well known to those of ordinary skill in the art. A basic text  
45 disclosing the general methods of use in this invention is Sambrook *et al.*, *Molecular Cloning, A Laboratory Manual* (2nd ed. 1989).  
30

Standard transfection methods are used to produce prokaryotic,  
50 mammalian, yeast or insect cell lines which express large quantities of the desired

5 polypeptide, which is then purified using standard techniques (*see, e.g., Colley et al., J. Biol. Chem.* 264:17619-17622, 1989; *Guide to Protein Purification, supra*).

10 The nucleotide sequences used to transfect the host cells can be modified to yield *Chlamydia* polypeptides with a variety of desired properties. For example, the polypeptides can vary from the naturally-occurring sequence at the primary structure level by amino acid, insertions, substitutions, deletions, and the like. These modifications can be used in a number of combinations to produce the final modified protein chain.

15 The amino acid sequence variants can be prepared with various objectives in mind, including facilitating purification and preparation of the recombinant polypeptide. The modified polypeptides are also useful for modifying plasma half life, improving therapeutic efficacy, and lessening the severity or occurrence of side effects during therapeutic use. The amino acid sequence variants are usually predetermined variants not found in nature but exhibit the same immunogenic activity as naturally occurring protein. In general, modifications of the sequences encoding the polypeptides may be readily accomplished by a variety of well-known techniques, such as site-directed mutagenesis (*see Gillman & Smith, Gene* 8:81-97 (1979); *Roberts et al., Nature* 328:731-734 (1987)). One of ordinary skill will appreciate that the effect of many mutations is difficult to predict. Thus, most modifications are evaluated by routine screening in a suitable assay for the desired characteristic. For instance, the effect of various modifications on the ability of the polypeptide to elicit a protective immune response can be easily determined using *in vitro* assays. For instance, the polypeptides can be tested for their ability to induce lymphoproliferation, T cell cytotoxicity, or cytokine production using standard techniques.

25 The particular procedure used to introduce the genetic material into the host cell for expression of the polypeptide is not particularly critical. Any of the well known procedures for introducing foreign nucleotide sequences into host cells may be used. These include the use of calcium phosphate transfection, spheroplasts, electroporation, liposomes, microinjection, plasmid vectors, viral vectors and any of the other well known methods for introducing cloned genomic DNA, cDNA, synthetic DNA or other foreign genetic material into a host cell (*see Sambrook et al., supra*). It is only necessary that the particular procedure utilized be capable of successfully introducing at least one gene into the host cell which is capable of expressing the gene.



5 Any of a number of well known cells and cell lines can be used to express the polypeptides of the invention. For instance, prokaryotic cells such as *E. coli* can be used. Eukaryotic cells include, yeast, Chinese hamster ovary (CHO) cells, COS cells, and insect cells.

10 5 The particular vector used to transport the genetic information into the cell is also not particularly critical. Any of the conventional vectors used for expression of recombinant proteins in prokaryotic and eukaryotic cells may be used. Expression vectors for mammalian cells typically contain regulatory elements from eukaryotic viruses.

15 10 The expression vector typically contains a transcription unit or expression cassette that contains all the elements required for the expression of the polypeptide DNA in the host cells. A typical expression cassette contains a promoter operably linked to the DNA sequence encoding a polypeptide and signals required for efficient polyadenylation of the transcript. The term "operably linked" as used herein refers to linkage of a promoter upstream from a DNA sequence such that the promoter mediates transcription of the DNA sequence. The promoter is preferably positioned about the same distance from the heterologous transcription start site as it is from the transcription start site in its natural setting. As is known in the art, however, some variation in this distance can be accommodated without loss of promoter function.

20 20 Following the growth of the recombinant cells and expression of the polypeptide, the culture medium is harvested for purification of the secreted protein. The media are typically clarified by centrifugation or filtration to remove cells and cell debris and the proteins are concentrated by adsorption to any suitable resin or by use of ammonium sulfate fractionation, polyethylene glycol precipitation, or by ultrafiltration. 25 Other routine means known in the art may be equally suitable. Further purification of the polypeptide can be accomplished by standard techniques, for example, affinity chromatography, ion exchange chromatography, sizing chromatography, His<sub>6</sub> tagging and Ni-agarose chromatography (as described in Dobeli *et al.*, *Mol. and Biochem. Parasit.* 41:259-268 (1990)), or other protein purification techniques to obtain homogeneity. The 40 purified proteins are then used to produce pharmaceutical compositions, as described below.

45 30 An alternative method of preparing recombinant polypeptides useful as vaccines involves the use of recombinant viruses (e.g., vaccinia). Vaccinia virus is grown

in suitable cultured mammalian cells such as the HeLa S3 spinner cells, as described by Mackett *et al.*, in *DNA cloning Vol. II: A practical approach*, pp. 191-211 (Glover, ed.).

#### Antibody Production

The proteins of the present invention can be used to produce antibodies specifically reactive with *C pneumoniae* antigens. If isolated proteins are used, they may be recombinantly produced or isolated from *Chlamydia* cultures. Synthetic peptides made using the protein sequences may also be used.

Methods of production of polyclonal antibodies are known to those of skill in the art. In brief, an immunogen, preferably a purified protein, is mixed with an adjuvant and animals are immunized. When appropriately high titers of antibody to the immunogen are obtained, blood is collected from the animal and antisera is prepared. Further fractionation of the antisera to enrich for antibodies reactive to *Chlamydia* proteins can be done if desired (*see Harlow & Lane, Antibodies: A Laboratory Manual* (1988)).

Polyclonal antisera are used to identify and characterize *Chlamydia* in the tissues of patients using, for instance, *in situ* techniques and immunoperoxidase test procedures described in Anderson *et al. JAVMA* 198:241 (1991) and Barr *et al. Vet. Pathol.* 28:110-116 (1991).

Monoclonal antibodies may be obtained by various techniques familiar to those skilled in the art. Briefly, spleen cells from an animal immunized with a desired antigen are immortalized, commonly by fusion with a myeloma cell (*see Kohler & Milstein, Eur. J. Immunol.* 6:511-519 (1976)). Alternative methods of immortalization include transformation with Epstein Barr Virus, oncogenes, or retroviruses, or other methods well known in the art. Colonies arising from single immortalized cells are screened for production of antibodies of the desired specificity and affinity for the antigen, and yield of the monoclonal antibodies produced by such cells may be enhanced by various techniques, including injection into the peritoneal cavity of a vertebrate host.

Monoclonal antibodies produced in such a manner are used, for instance, in ELISA diagnostic tests, immunoperoxidase tests, immunohistochemical tests, for the *in vitro* evaluation of spirochete invasion, to select candidate antigens for vaccine development, protein isolation, and for screening genomic and cDNA libraries to select appropriate gene sequences.

Immunodiagnostic detection of *C. pneumoniae* infections

The present invention also provides methods for detecting the presence or absence of *C. pneumoniae*, or antibodies reactive with it, in a biological sample. For instance, antibodies specifically reactive with *Chlamydia* can be detected using either *Chlamydia* proteins or the isolates described here. The proteins and isolates can also be used to raise specific antibodies (either monoclonal or polyclonal) to detect the antigen in a sample. In addition, the nucleic acids disclosed and claimed here can be used to detect *Chlamydia*-specific sequences using standard hybridization techniques.

For a review of immunological and immunoassay procedures in general, see *Basic and Clinical Immunology* (Stites & Terr ed., 7th ed. 1991)). The immunoassays of the present invention can be performed in any of several configurations, which are reviewed extensively in *Enzyme Immunoassay* (Maggio, ed., 1980); Tijssen, *Laboratory Techniques in Biochemistry and Molecular Biology* (1985)). For instance, the proteins and antibodies disclosed here are conveniently used in ELISA, immunoblot analysis and agglutination assays.

In brief, immunoassays to measure anti-*Chlamydia* antibodies or antigens can be either competitive or noncompetitive binding assays. In competitive binding assays, the sample analyte (e.g., anti-*Chlamydia* antibodies) competes with a labeled analyte (e.g., anti-*Chlamydia* monoclonal antibody) for specific binding sites on a capture agent (e.g., isolated *Chlamydia* protein) bound to a solid surface. The concentration of labeled analyte bound to the capture agent is inversely proportional to the amount of free analyte present in the sample.

Noncompetitive assays are typically sandwich assays, in which the sample analyte is bound between two analyte-specific binding reagents. One of the binding agents is used as a capture agent and is bound to a solid surface. The second binding agent is labelled and is used to measure or detect the resultant complex by visual or instrument means.

A number of combinations of capture agent and labelled binding agent can be used. For instance, an isolated *Chlamydia* protein or culture can be used as the capture agent and labelled anti-human antibodies specific for the constant region of human antibodies can be used as the labelled binding agent. Goat, sheep and other non-human antibodies specific for human immunoglobulin constant regions (e.g.,  $\gamma$  or  $\mu$ ) are

5 well known in the art. Alternatively, the anti-human antibodies can be the capture agent and the antigen can be labelled.

10 Various components of the assay, including the antigen, anti-*Chlamydia* antibody, or anti-human antibody, may be bound to a solid surface. Many methods for  
5 immobilizing biomolecules to a variety of solid surfaces are known in the art. For instance, the solid surface may be a membrane (e.g., nitrocellulose), a microtiter dish (e.g., PVC or polystyrene) or a bead. The desired component may be covalently bound or  
15 noncovalently attached through nonspecific bonding.

10 Alternatively, the immunoassay may be carried out in liquid phase and a variety of separation methods may be employed to separate the bound labeled component from the unbound labelled components. These methods are known to those of skill in the  
20 art and include immunoprecipitation, column chromatography, adsorption, addition of magnetizable particles coated with a binding agent and other similar procedures.

25 An immunoassay may also be carried out in liquid phase without a separation procedure. Various homogeneous immunoassay methods are now being applied to immunoassays for protein analytes. In these methods, the binding of the binding agent to the analyte causes a change in the signal emitted by the label, so that  
30 binding may be measured without separating the bound from the unbound labelled component.

20 Western blot (immunoblot) analysis can also be used to detect the presence of antibodies to *Chlamydia* in the sample. This technique is a reliable method for  
35 confirming the presence of antibodies against a particular protein in the sample. The technique generally comprises separating proteins by gel electrophoresis on the basis of molecular weight, transferring the separated proteins to a suitable solid support, (such as a  
25 nitrocellulose filter, a nylon filter, or derivatized nylon filter), and incubating the sample with the separated proteins. This causes specific target antibodies present in the sample  
40 to bind their respective proteins. Target antibodies are then detected using labeled anti-human antibodies.

45 The immunoassay formats described above employ labelled assay components. The label may be coupled directly or indirectly to the desired component of  
30 the assay according to methods well known in the art. A wide variety of labels may be used. The component may be labelled by any one of several methods. Traditionally a  
50 radioactive label incorporating  $^3\text{H}$ ,  $^{125}\text{I}$ ,  $^{35}\text{S}$ ,  $^{14}\text{C}$ , or  $^{32}\text{P}$  was used. Non-radioactive labels

5 include ligands which bind to labelled antibodies, fluorophores, chemiluminescent agents, enzymes, and antibodies which can serve as specific binding pair members for a labelled ligand. The choice of label depends on sensitivity required, ease of conjugation with the compound, stability requirements, and available instrumentation.

10 5 Enzymes of interest as labels will primarily be hydrolases, particularly phosphatases, esterases and glycosidases, or oxidoreductases, particularly peroxidases. Fluorescent compounds include fluorescein and its derivatives, rhodamine and its derivatives, dansyl, umbelliferone, etc. Chemiluminescent compounds include luciferin, and 2,3-dihydrophthalazinediones, e.g., luminol. For a review of various labelling or  
15 signal producing systems which may be used, see U.S. Patent No. 4,391,904, which is incorporated herein by reference.

20 Non-radioactive labels are often attached by indirect means. Generally, a ligand molecule (e.g., biotin) is covalently bound to the molecule. The ligand then binds to an anti-ligand (e.g., streptavidin) molecule which is either inherently detectable or  
25 covalently bound to a signal system, such as a detectable enzyme, a fluorescent compound, or a chemiluminescent compound. A number of ligands and anti-ligands can be used. Where a ligand has a natural anti-ligand, for example, biotin, thyroxine, and cortisol, it can be used in conjunction with the labelled, naturally occurring anti-ligands. Alternatively, any haptenic or antigenic compound can be used in combination with an  
30 antibody.

35 Some assay formats do not require the use of labelled components. For instance, agglutination assays can be used to detect the presence of the target antibodies. In this case, antigen-coated particles are agglutinated by samples comprising the target antibodies. In this format, none of the components need be labelled and the presence of  
25 the target antibody is detected by simple visual inspection.

#### 40 Pharmaceutical Compositions

45 The peptides or antibodies (typically monoclonal antibodies) of the present invention and pharmaceutical compositions thereof are useful for administration to mammals, particularly humans, to treat and/or prevent *Chlamydia* infections. Suitable  
30 formulations are found in *Remington's Pharmaceutical Sciences*, Mack Publishing Company, Philadelphia, PA, 17th ed. (1985).

5                   The immunogenic peptides or antibodies of the invention are administered  
prophylactically or to an individual already suffering from the disease. The peptide  
compositions are administered to a patient in an amount sufficient to elicit an effective  
10                   immune response to *Chlamydia*. An effective immune response is one that inhibits  
5                   infection. An amount adequate to accomplish this is defined as "therapeutically effective  
dose" or "immunogenically effective dose." Amounts effective for this use will depend  
on, e.g., the peptide composition, the manner of administration, the stage and severity of  
15                   the disease being treated, the weight and general state of health of the patient, and the  
judgment of the prescribing physician, but generally range for the initial immunization  
10                   (that is for therapeutic or prophylactic administration) from about 0.1 mg to about 1.0 mg  
per 70 kilogram patient, more commonly from about 0.5 mg to about 0.75 mg per 70 kg  
20                   of body weight. Boosting dosages are typically from about 0.1 mg to about 0.5 mg of  
peptide using a boosting regimen over weeks to months depending upon the patient's  
response and condition. A suitable protocol would include injection at time 0, 4, 2, 6, 10  
25                   15                   and 14 weeks, followed by further booster injections at 24 and 28 weeks.

For therapeutic use, administration should begin at the first sign of  
infection. This is followed by boosting doses until at least symptoms are substantially  
abated and for a period thereafter. In some circumstances, loading doses followed by  
30                   boosting doses may be required. The resulting immune response helps to cure or at least  
20                   partially arrest symptoms and/or complications. Vaccine compositions containing the  
peptides are administered prophylactically to a patient susceptible to or otherwise at risk  
of the infection.

35                   The pharmaceutical compositions (containing either peptides or  
antibodies) are intended for parenteral or oral administration. Preferably, the  
25                   pharmaceutical compositions are administered parenterally, e.g., subcutaneously,  
40                   intradermally, or intramuscularly. Thus, the invention provides compositions for  
parenteral administration which comprise a solution of the immunogenic polypeptides  
dissolved or suspended in an acceptable carrier, preferably an aqueous carrier. A variety  
of aqueous carriers may be used, e.g., water, buffered water, 0.4% saline, 0.3% glycine,  
45                   30                   hyaluronic acid and the like. These compositions may be sterilized by conventional, well  
known sterilization techniques, or may be sterile filtered. The resulting aqueous solutions  
may be packaged for use as is, or lyophilized, the lyophilized preparation being combined  
50                   with a sterile solution prior to administration. The compositions may contain

5 pharmaceutically acceptable auxiliary substances as required to approximate physiological conditions, such as buffering agents, tonicity adjusting agents, wetting agents and the like, for example, sodium acetate, sodium lactate, sodium chloride, potassium chloride, calcium chloride, sorbitan monolaurate, triethanolamine oleate, etc.

10 5 The compositions may also comprise carriers to enhance the immune response. Useful carriers are well known in the art, and include, e.g., KLH, thyroglobulin, albumins such as human serum albumin, tetanus toxoid, polyamino acids such as poly(lysine:glutamic acid), influenza, hepatitis B virus core protein, hepatitis B virus recombinant vaccine and the like.

10 For solid compositions, conventional nontoxic solid carriers may be used which include, for example, pharmaceutical grades of mannitol, lactose, starch, magnesium stearate, sodium saccharin, talcum, cellulose, glucose, sucrose, magnesium carbonate, and the like. For oral administration, a pharmaceutically acceptable nontoxic composition is formed by incorporating any of the normally employed excipients, such as  
20 those carriers previously listed, and generally 10-95% of active ingredient, that is, one or more peptides of the invention, and more preferably at a concentration of 25%-75%.

As noted above, the peptide compositions are intended to induce an immune response to *Chlamydia*. Thus, compositions and methods of administration suitable for maximizing the immune response are preferred. For instance, peptides may  
20 be introduced into a host, including humans, linked to a carrier or as a homopolymer or heteropolymer of active peptide units from various *Chlamydia* proteins disclosed here. Alternatively, a "cocktail" of polypeptides can be used. A mixture of more than one polypeptide has the advantage of increased immunological reaction and, where different peptides are used to make up the polymer, the additional ability to induce antibodies to a  
25 number of epitopes.

40 The compositions also include an adjuvant. As used here, number of adjuvants are well known to one skilled in the art. Suitable adjuvants include incomplete Freund's adjuvant, alum, aluminum phosphate, aluminum hydroxide, N-acetyl-muramyl-L-threonyl-D-isoglutamine (thr-MDP),  
45 N-acetyl-nor-muramyl-L-alanyl-D-isoglutamine (CGP 11637, referred to as nor-MDP), N-acetylmuramyl-L-alanyl-D-isoglutaminyl-L-alanine-2-(1'-2'-dipalmitoyl-sn-glycero-3-hydroxyphosphoryloxy)-ethylamine (CGP 19835A, referred to as MTP-PE),  
50 and RIBI, which contains three components extracted from bacteria, monophosphoryl

lipid A, trehalose dimycolate and cell wall skeleton (MPL+TDM+CWS) in a 2% squalene/Tween 80 emulsion. The effectiveness of an adjuvant may be determined by measuring the amount of antibodies directed against the immunogenic peptide.

The concentration of immunogenic peptides of the invention in the pharmaceutical formulations can vary widely, i.e. from less than about 0.1%, usually at or at least about 2% to as much as 20% to 50% or more by weight, and will be selected primarily by fluid volumes, viscosities, etc., in accordance with the particular mode of administration selected.

The peptides of the invention can also be expressed by attenuated viral hosts, such as vaccinia or fowlpox. This approach involves the use of vaccinia virus as a vector to express nucleotide sequences that encode the peptides of the invention. Upon introduction into a host, the recombinant vaccinia virus expresses the immunogenic peptide, and thereby elicits an immune response. Vaccinia vectors and methods useful in immunization protocols are described in, e.g., U.S. Patent No. 4,722,848. Another vector is BCG (Bacille Calmette Guerin). BCG vectors are described in Stover et al. (*Nature* 351:456-460 (1991)). A wide variety of other vectors useful for therapeutic administration or immunization of the peptides of the invention, e.g., *Salmonella typhi* vectors and the like, will be apparent to those skilled in the art from the description herein.

The DNA encoding one or more of the peptides of the invention can also be administered to the patient. This approach is described, for instance, in Wolff *et al.*, *Science* 247: 1465-1468 (1990) as well as U.S. Patent Nos. 5,580,859 and 5,589,466.

In order to enhance serum half-life, the peptides may also be encapsulated, introduced into the lumen of liposomes, prepared as a colloid, or other conventional techniques may be employed which provide an extended serum half-life of the peptides. A variety of methods are available for preparing liposomes, as described in, e.g., Szoka et al., *Ann. Rev. Biophys. Bioeng.* 9:467 (1980), U.S. Pat. Nos. 4,235,871, 4,501,728 and 4,837,028.

### EXAMPLES

The following examples are offered to illustrate, but no to limit the claimed invention.

#### Example 1:



5 This example describes comparison of the *C. pneumoniae* genome disclosed here and the, previously sequenced, *C. trachomatis* genome (Stephens, *et al. Science* 282:754-759 (1998)).

10 The apparent low level of DNA homology between *C. trachomatis* and *C. pneumoniae* (Campbell, *et al., J. Clin. Microbiol.* 25:1911-1916 (1987)) yet analogous cell structures and developmental cycles, predicts that comparative analysis of the two genomes will significantly enhance the understanding of both pathogens. Identification of genes that are present in one species but not the other are of particular importance for the mutually exclusive biological, virulence and pathogenesis capabilities of each.

15 Identification of genes shared between the two species strongly supports the requirement for these capabilities in a biological system that has, over its long-term association with mammalian host cells, evolved to reduce the metabolic capacities while optimizing survival, growth and transmission of these unique pathogens.

20 The previously sequenced *C. trachomatis* genome contains 1,042,519 nucleotides and 875 likely protein-coding genes. Similarity searching permitted the inferred functional assignment of sequences 636 (60%) genes disclosed here and 251 (23%) are similar to hypothetical genes for other bacterial organisms including those for *C. trachomatis*. The remaining 186 (17%) genes are not homologous to sequences deposited in GenBank.. Seventy *C. trachomatis* genes are not represented in the *C. pneumoniae* genome. These are contained within blocks consisting of 2-17 genes and 19 single genes. Of the 70 *C. trachomatis* genes without homologs in *C. pneumoniae*, 60 are classified as encoding hypothetical proteins. The remaining genes not represented in *C. pneumoniae* consist of the tryptophan operon (*trpA,B,R*), *trpC*, two predicted thiol protease genes, and 4 genes assigned to the phospholipase-D superfamily.

25 It is evident that there is a high level of functional conservation between *C. pneumoniae* and *C. trachomatis* as orthologs to *C. trachomatis* genes were identified for 859 (80%) of the predicted coding sequences for *C. pneumoniae*. The level of similarity for individual encoded proteins spans a wide spectrum (22-95% amino acid identity) with an average of 62% amino acid identity between orthologs from the two species. The percent amino acid identity between orthologous chlamydial proteins is similar among functional groups with the highest for proteins associated with translation and the lowest for proteins whose function in chlamydiae is uncharacterized and not related to proteins encoded by other organisms. The gene order of the homologous set of genes in *C.*

*pneumoniae* shows reorganization relative to the genome of *C. trachomatis*; however, there is a high level of synteny for the gene organization of the two genomes. We identified thirty-nine blocks of 2 or more genes whose gene organization is colinear with homologs to *C. trachomatis*, although some of these are inverted. The distribution of genome reorganization is not evenly distributed on the chromosome as the region between *C. pneumoniae* coding sequences 0130-0300 contains substantially more reorganization than other areas of the genome. This region coincides with the predicted chromosome replication terminus.

We identified orthologs of enzymes characterized in other bacteria that account for the essential requirements for DNA replication, repair, transcription and translation including two predicted DNA helicases of the Swi2/Snf2 family found in *C. trachomatis*. Similar to *C. trachomatis*, alternative sigma subunits for RNA polymerase,  $\sigma^{28}$  and  $\sigma^{54}$ , were identified in addition to anti- $\sigma$  regulatory system factors RsbV, a RsbW-like single-domain histidine kinase, and a RsbU-like protein phosphatase. These findings suggest that the fundamental mechanisms of transcriptional regulation are conserved among *Chlamydia*. The *C. trachomatis* proteins containing SET and SWIB domains, and a SWIB domain fused to the C-terminus of the chlamydial topoisomerase I, not identified outside eukaryotes, are found in *C. pneumoniae* supporting their possible role in the chromatin condensation-decondensation characteristic of the biologically unique chlamydial developmental cycle.

The central metabolic pathways inferred from the *C. pneumoniae* genome sequence are the same as those identified for *C. trachomatis*. *C. pneumoniae* has a glycolytic pathway and a linked tricarboxylic acid cycle, although likely functional, is incomplete as genes for citrate synthase, aconitase, and isocitrate dehydrogenase were not identified. *C. pneumoniae* has a complete glycogen synthesis and degradation system supporting a role for glycogen synthesis and utilization of glucose-derivatives in chlamydial metabolism. Genes encoding essential functions in aerobic respiration are present and electron flux may be supported by pyruvate, succinate, glycerol-3-phosphate, and NADH dehydrogenases, NADH-ubiquinone oxidoreductase and cytochrome oxidase. *C. pneumoniae* also contains the V (vacuolar)-type ATPase operon and the two ATP translocases found in *C. trachomatis*.

The type-III secretion virulence system required for invasion by several pathogenic bacteria and found in the *C. trachomatis* genome in three chromosomal

locations is also present in the *C. pneumoniae* genome. Each of the components is conserved and their relative genomic contexts are conserved. Genes such as a predicted serine/threonine protein kinase and other genes physically linked to genes encoding structural components of the type-III secretion apparatus, but without identified homologs, are also highly similar between the two species suggesting the functional roles in modifying cellular biology are fundamentally conserved.

*Chlamydia*-encoded proteins that are not found in chlamydial organisms but localized to the intracellular chlamydial inclusion membrane are likely essential for the unique intracellular biology and perhaps differences in inclusion morphology observed between species of *Chlamydia*. Several such proteins, termed IncA, B&C, have been characterized for a *C. psittaci* strain (Rockey, et al. *Mol. Microbiol.* 15:617-626 (1995); Rockey et al. *Infect. Immun.* 62:106-112 (1994)). *C. pneumoniae* and *C. trachomatis* encode orthologs to *C. psittaci* IncB and IncC and *C. trachomatis* also contains an ortholog to IncA. *C. pneumoniae* contains two genes that encode proteins with similarity to IncA (CPn0186 and CPn0585), although the level of homology is low suggesting analogous but possibly altered functions.

The tryptophan biosynthesis operon (*trpA*, *trpB*, *trpR*) and *trpC* identified in *C. trachomatis* is conspicuously missing in the *C. pneumoniae* genome. This represents the entire repertoire of genes associated with tryptophan biosynthesis identified in *C. trachomatis*. Seventeen genes adjacent to the *C. trachomatis* tryptophan operon also were not found in the *C. pneumoniae* genome. This region is the single largest loss of a contiguous genomic segment and includes 4 HKD superfamily encoding genes that encompass a family of proteins related to endonuclease and phospholipase D. These findings may be important for the ability of *Chlamydia* to persist in their hosts and cause disease by eliciting potent, focal and persistent inflammatory responses thought to be essential for pathogenesis.

The *C. pneumoniae* genome contains 187,711 additional nucleotides compared to the *C. trachomatis* genome, and the 214 coding sequences not found in *C. trachomatis* account for most of the increased genome size. Eighty-eight of these genes are found in blocks of >10 genes (11-30 genes/block), 41 are single genes, and the remainder are partnered with at least one other gene. Based upon the observation that ~70% of all the *C. pneumoniae* genes have an identifiable homolog in GenBank, exclusive of *C. trachomatis*, it would be expected that over 150 of the 214 genes should

5 have a homolog in GenBank, many associated with a function. However, only 28 coding  
sequences have similarity to genes from other organisms. Thus the majority of the genes  
that are mutually exclusive of *C. trachomatis* (186 of 214), and the 60 of 70 *C.*  
10 *trachomatis* genes that lacked an identifiable homolog in *C. pneumoniae*, do not have  
5 detectable homologs to genes from other organisms. We predict that most of the unique  
genes are essential for specific attributes that define the differential biology, tropism and  
pathogenesis of *C. trachomatis* and *C. pneumoniae*. Moreover, this suggests that *C.*  
15 *pneumoniae* has more unique biological (i.e., virulence) capacity than *C. trachomatis*.  
The ability of *C. pneumoniae* to be more invasive and survive in a broader range of host  
10 cell types than *C. trachomatis* is consistent with this hypothesis. Not all of the  
differences in biological capacity may be associated with mutually exclusive genes. One  
20 explanation for the significantly lower level of homology between protein sequences  
assigned as having *C. pneumoniae* and *C. trachomatis* orthologs but no identifiable  
orthologs in other organisms is that this set of proteins is not only associated with  
25 biological requirements specific for *Chlamydia* but this polymorphism may account for  
differential biology between the two species. The determination of the genome sequence  
from a representative of the *C. psittaci* group will precisely delineate those genes that are  
mutually exclusive and specific for each species.

30 The major functionally identifiable addition to the *C. pneumoniae* genome  
20 is a large expansion of genes encoding a new family of chlamydial polymorphic  
membrane proteins (Pmp), alone representing 22% of the increased coding capacity.  
While the *C. trachomatis* genome has 9 *pmp* genes, remarkably the *C. pneumoniae*  
35 genome contains 21 *pmp* genes. Most of these genes appear to be amplified in two  
regions of the genome with three stand-alone genes. Interestingly one of the stand-alone  
25 genes is most closely related to the *C. trachomatis pmpD* which is the only stand-alone  
*pmp* gene in the *C. trachomatis* genome and it is located with the same relative genomic  
40 context, suggesting an essential and conserved function for this paralog. Six Pmp-coding  
genes are presumably not functional as five contain predicted coding frame-shifts and one  
is truncated. The amplification of this gene family and the confidently predicted frame-  
45 shifts suggest a specific molecular mechanism to promote functional or antigenic  
30 diversity. The biological role of this protein family remains enigmatic, although at least  
one of the proteins in *C. psittaci* related to this family is exposed on the chlamydial  
50 surface.

5 While a function could not be assigned for most of the unique *C. pneumoniae* genes, several have significant similarity to genes from other organisms. Functional assignments could be made for genes encoding GMP synthetase, IMP dehydrogenase, UMP synthase, uridine kinase, biotin synthase pathway proteins,  
10 5 methylthioadenosine nucleosidase, a DNA glycosylase and aromatic amino acid hydroxylase. Thus a complete pathway was identified for biotin biosynthesis. The additional purine and pyrimidine salvage pathway genes presumably reflect metabolic limitations in one of the cell types that *C. pneumoniae* infects or differences in the ability  
15 of *C. pneumoniae* to transport precursor nucleosides or nucleotides.

10 The addition of aromatic amino acid hydroxylase in *C. pneumoniae* is intriguing especially in light of the loss of tryptophan biosynthetic genes and the inability to synthesize other amino acids including phenylalanine. Aromatic amino acid  
20 hydroxylases include three distinct enzymes that function to receptively oxidize phenylalanine to tyrosine, tyrosine to Dopa, and tryptophan to 5-hydroxytryptophan and serotonin. Although the chlamydial protein is similar to proteins of this family and  
25 15 incrementally more closely related to tryptophan hydroxylase, its specific function could not be confidently predicted. We hypothesize that it may be involved in *C. pneumoniae* virulence. Tryptophan hydroxylase has not been previously identified in bacteria and the origin of the chlamydial gene appears to be from eukaryotes. The functional role of an  
30 20 aromatic amino acid hydroxylase for *C. pneumoniae* is linked to the unique intracellular biology of this organism and may represent a key contribution to *C. pneumoniae* persistence and pathogenesis.

35 It is understood that the examples and embodiments described herein are for illustrative purposes only and that various modifications or changes in light thereof  
25 will be suggested to persons skilled in the art and are to be included within the spirit and purview of this application and scope of the appended claims. All publications, patents,  
40 and patent applications cited herein are hereby incorporated by reference in their entirety for all purposes.

45 Table 1 provides functional assignments of *C. pneumoniae* nonprotein-  
30 encoding genomic sequences. Table 2 provides functional assignments of protein coding sequences. Table 3 provides the amino acid sequences of the proteins corresponding to the coding sequences.

5

What is Claimed is:

1. An isolated nucleic acid encoding a *C. pneumoniae* protein as set forth in Table 3.

10

5

2. The isolated nucleic acid of Claim 1, wherein said nucleic acid has a nucleotide sequence of an open reading frame in SEQ ID NO:1.

15

3. A probe comprising a hybridizing fragment of an isolated nucleic acid according to Claim 2.

10

20

5. An isolated nucleic acid that hybridizes under stringent conditions to the nucleic acid sequence of Claim 2.

15

6. An expression cassette comprising a transcriptional initiation region functional in an expression host, a nucleic acid having a sequence of the isolated nucleic acid according to Claim 1 under the transcriptional regulation of said transcriptional initiation region, and a transcriptional termination region functional in said expression host.

25

30

20

7. A cell comprising an expression cassette according to Claim 6 as part of an extrachromosomal element or integrated into the genome of a host cell as a result of introduction of said expression cassette into said host cell, and the cellular progeny of said host cell.

35

25

8. A method for producing a *C. pneumoniae* protein, said method comprising:  
growing a cell according to Claim 7, whereby said *C. pneumoniae* protein is expressed; and  
isolating said *C. pneumoniae* protein free of other proteins.

40

45

30

50

55

5

9. A purified polypeptide composition comprising at least 50 weight % of the protein present as a *C. pneumoniae* protein comprising an amino acid sequence of claim 1.

10

5

10. A monoclonal antibody binding specifically to the polypeptide of Claim 9.

15

20

25

30

35

40

45

50

55

**THIS PAGE BLANK (USPTO)**